

Week 9 Notes: Fair Prediction with Disparate Impact

Tuesday October 22, 2019

Notation

- ❖ Let $S=s(x)$ denote the risk Score.
- ❖ We let $R \in \{b, w\}$ denote the group to which an individual belongs
- ❖ We denote the outcome indicator by $Y \in \{0,1\}$, with $Y=1$ indicating that the given individual goes on to recidivate.
- ❖ We introduce the quantity s_{HR} , which denotes the high-risk score threshold.
- ❖ Defendants whose score S exceeds s_{HR} will be referred to as high risk, while the remaining defendants will be referred to as low risk.

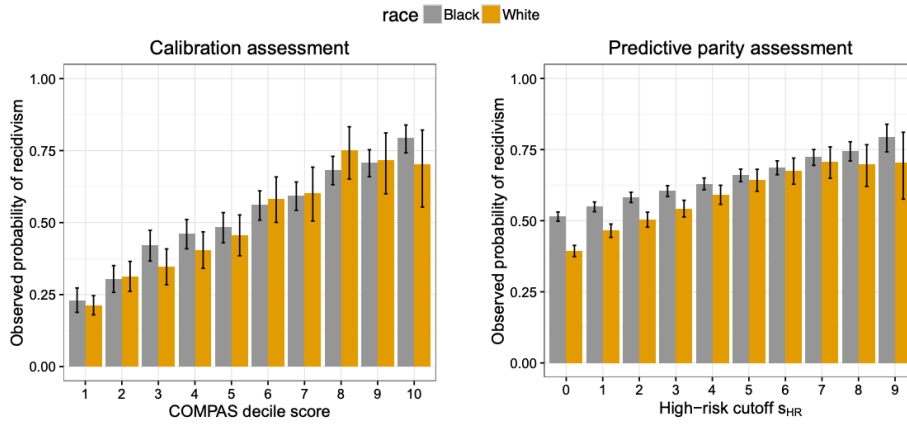
Several Different Fairness Criteria

- ❖ Calibration
 - A score $s=s(x)$ is said to be well-calibrated if it reflects the same likelihood of recidivism irrespective of the individual's group membership.
 - For all values of s
 - $\mathbb{P}(Y = 1 | S = s, R = b) = \mathbb{P}(Y = 1 | S = s, R = w)$.
 - In their response to the ProPublica investigation, Flores et. Al [6] verify that COMPAS is well-calibrated using logistic regression modeling.
- ❖ Predictive Parity
 - A score of $S=S(x)$ satisfies predictive parity at a threshold s_{HR} if the likelihood of recidivism among high-risk offenders is the same regardless of the group membership.
 - $\mathbb{P}(Y = 1 | S > s_{HR}, R = b) = \mathbb{P}(Y = 1 | S > s_{HR}, R = w)$.
 - Predictive parity at a given threshold s_{HR} amounts to requiring that the positive predictive value (PPV) of the classifier $Y^{\wedge} = 1_{S>s_{HR}}$ be the same across groups.
 - Northpointe's refutation of the ProPublica analysis shows that COMPAS satisfies predictive parity for threshold choices of interest.
- ❖ Error Rate Balance
 - A score $S=S(x)$ satisfies error rate balance at a threshold s_{HR} if the false positive and false negative error rates are equal across groups.
 - $\mathbb{P}(S > s_{HR}, | Y = 0, R = b) = \mathbb{P}(S > s_{HR}, Y = 0, R = w)$
 - $\mathbb{P}(S < s_{HR}, | Y = 1, R = b) = \mathbb{P}(S < s_{HR}, Y = 1, R = w)$
 - ProPublica's analysis considered a threshold of $s_{HR} = 4$, which they showed leads to considerable imbalance in both false positive and false negative rates/
 - Error rate balance is also closely connected to the notions of equalized odds and equal opportunity.
- ❖ Statistical Parity
 - A score $S=S(x)$ satisfies statistical parity at a threshold s_{HR} if the proportion of individuals classified as high-risk is the same for each group.
 - $\mathbb{P}(S > s_{HR}, | R = b) = \mathbb{P}(S > s_{HR}, | R = w)$

- Other names: Demographic parity, equal acceptance rates, group fairness.

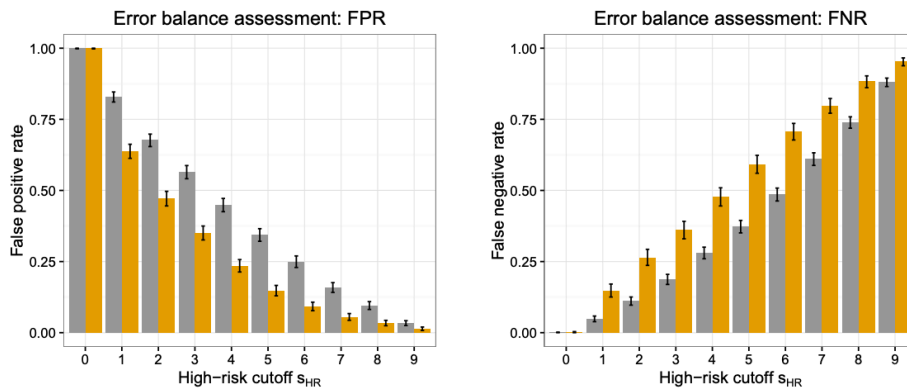
What Does COMPAS Satisfy?

- ❖ Calibration Assessment
 - Yes
 - 95% Confidence Intervals Intercept for Black and White groups
- ❖ Predictive Parity
 - Yes
 - 95% Confidence Intervals Intercept for Black and White groups
- ❖ Error Balance
 - No
 - Higher False Positive Rates for Black Groups is Higher
- ❖ Statistical Parity
 - No
 - Higher False Negative Rates for White Groups is Higher



(a) Bars represent empirical estimates of the expressions in (2.1): $\mathbb{P}(Y = 1 \mid S = s, R = r)$ for decile scores $s \in \{1, \dots, 10\}$.

(b) Bars represent empirical estimates of the expressions in (2.2): $\mathbb{P}(Y = 1 \mid S > s_{HR}, R = r)$ for values of the high-risk cutoff $s_{HR} \in \{0, \dots, 9\}$.



(c) Bars represent observed false positive rates, which are empirical estimates of the expressions in (2.3): $\mathbb{P}(S > s_{HR} \mid Y = 0, R = r)$ for values of the high-risk cutoff $s_{HR} \in \{0, \dots, 9\}$.

(d) Bars represent observed false negative rates, which are empirical estimates of the expressions in (2.4): $\mathbb{P}(S \leq s_{HR} \mid Y = 1, R = r)$ for values of the high-risk cutoff $s_{HR} \in \{0, \dots, 9\}$.

What is the Main Finding?

- ❖ The error rate imbalance exhibited by COMPAS is not a coincidence, nor can it be remedied in the present context
- ❖ Impossibility Result
- ❖ **When the recidivism prevalence i.e., the base rate $P(Y=1 | R=r)$ differs across groups, any instrument that satisfies predictive parity (PPR) at a given threshold sHR must have imbalanced false positive or false negative rates at that threshold.**

$$FPR = \frac{p}{1-p} \frac{1-PPV}{PPV} (1-FNR).$$

Thursday October 24, 2019

Does This lead to disparate impact?

- ❖ What is disparate impact?
- ❖ Definition could be context specific
- ❖ In COMPAS context, let's say that defendant receives higher penalty if he/she is adjudged high-risk and less penalty

Corollary 3.1 (Non-Recidivists). *Among individuals who do not recidivate, the difference in average penalty under the MinMax policy is*

$$\Delta = (t_{\max} - t_{\min})(FPR_b - FPR_w), \quad (3.2)$$

where FPR_r denotes the false positive rate among individuals in group $R = r$.

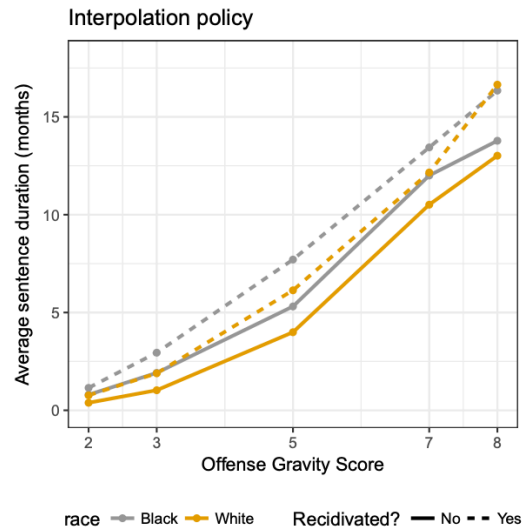
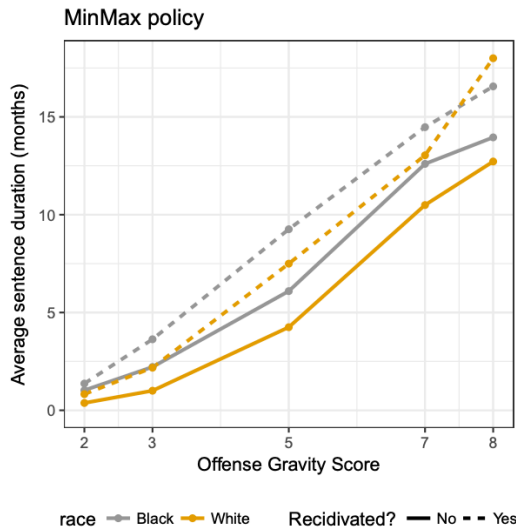
Corollary 3.2 (Recidivists). *Among individuals who recidivate, the difference in average penalty under the MinMax policy is*

$$\Delta = (t_{\max} - t_{\min})(FNR_w - FNR_b), \quad (3.3)$$

where FNR_r denotes the false negative rate among individuals in group $R = r$.

Is there disparate impact with COMPAS?

- ❖ When using an RPI that satisfies predictive parity in populations where recidivism prevalence differs across groups, it will generally be the case that the higher recidivism prevalence group will have higher FPR and lower FNR
 - In the COMPAS (RPI) setting
 - Does it satisfy predictive parity? Yes
 - Recidivism prevalence differs across groups? Yes
 - Which is the higher recidivism prevalence group? Black defendants
- ❖ This would on average result in greater penalties for defendants in the higher prevalence group, both among recidivists and non-recidivists.



Equality of Opportunity in Supervised Learning

Three Ways of Fixing Biases

- ❖ Preprocess Data
- ❖ Postprocessing Data
- ❖ Modify ML Algorithm

What is their high-level approach?

- ❖ Don't change anything in the ML training pipeline
 - Train models like you usually do
- ❖ Do some post-processing on the outputs of the ML model
 - Make it fairer
 - Make this process oblivious to the training set

2 Fairness criteria that they try to establish

- ❖ Equalize odds
 - Across your two categories
 - FPR should be equal
 - TPR should be equal

$$\Pr\{\widehat{Y} = 1 \mid A = 0, Y = y\} = \Pr\{\widehat{Y} = 1 \mid A = 1, Y = y\}, \quad y \in \{0, 1\}$$

- ❖ Equal Opportunity
- ❖ We say that a binary predictor Y_b satisfies equal opportunity with respect to A if Y if only satisfies for the positive class
- ❖ Equal opportunity is weaker, though still interesting, notion of non-discrimination

Deriving from Binary Predictor

$$\min_{\tilde{Y}} \mathbb{E} \ell(\tilde{Y}, Y) \quad (4.3)$$

$$\text{s.t. } \forall a \in \{0, 1\} : \gamma_a(\tilde{Y}) \in P_a(\hat{Y}) \quad (\text{derived})$$

$$\gamma_0(\tilde{Y}) = \gamma_1(\tilde{Y}) \quad (\text{equalized odds})$$

Deriving from Binary Predictor pt.2

$$\gamma_a(\hat{Y}) \stackrel{\text{def}}{=} \left(\Pr\{\hat{Y} = 1 \mid A = a, Y = 0\}, \Pr\{\hat{Y} = 1 \mid A = a, Y = 1\} \right). \quad (4.1)$$

Lemma 4.2. A predictor \hat{Y} satisfies:

1. equalized odds if and only if $\gamma_0(\hat{Y}) = \gamma_1(\hat{Y})$, and
2. equal opportunity if and only if $\gamma_0(\hat{Y})$ and $\gamma_1(\hat{Y})$ agree in the second component, i.e., $\gamma_0(\hat{Y})_2 = \gamma_1(\hat{Y})_2$.

$$P_a(\hat{Y}) \stackrel{\text{def}}{=} \text{convhull}\{(0, 0), \gamma_a(\hat{Y}), \gamma_a(1 - \hat{Y}), (1, 1)\} \quad (4.2)$$

Another Approach: Pre-Processing

- ❖ Remove the sensitive attribute
- ❖ Remove all features with sensitive attribute
- ❖ Brute Force Method?
 - For all the features that are somewhat correlated with the sensitive attribute
- ❖ What is a more sophisticated method?

